# 1    Integrative pathway enrichment analysis of multivariate omics data

2

3    Marta Paczkowska[1,*], Jonathan Barenboim[1,*], Nardnisa Sintupisut[1], Natalie C. Fox[1,2], Helen

4    Zhu[1,2], Diala Abd-Rabbo[1], PCAWG Network and Pathway Analysis Group, Paul C. Boutros[1,2,3],

5    Jüri Reimand[1,2,@]

6

7    1 - Ontario Institute for Cancer Research, Toronto, Ontario, Canada

8    2 - Department of Medical Biophysics, University of Toronto, Toronto, Ontario, Canada

9    3 - Department of Pharmacology & Toxicology, University of Toronto, Toronto, Ontario, Canada

10    * - these authors contributed equally

11    @ - correspondence: Juri.Reimand@utoronto.ca

12

13    **ABSTRACT**

14    **Multi-omics datasets quantify complementary aspects of molecular biology and thus pose**

15    **challenges to data interpretation and hypothesis generation. ActivePathways is an**

16    **integrative method that discovers significantly enriched pathways across multiple omics**

17    **datasets using a statistical data fusion approach, rationalizes contributing evidence and**

18    **highlights associated genes. We demonstrate its utility by analyzing coding and non-**

19    **coding mutations from 2,583 whole cancer genomes, revealing frequently mutated**

20    **hallmark pathways and a long tail of known and putative cancer driver genes. We also**

21    **studied prognostic molecular pathways in breast cancer subtypes by integrating genomic**

22    **and transcriptomic features of tumors and tumor-adjacent cells and found significant**

23    **associations with immune response processes and anti-apoptotic signaling pathways.**

24    **ActivePathways is a versatile method that improves systems-level understanding of**

25    **cellular organization in health and disease through integration of multiple molecular**

26    **datasets and pathway annotations.**

**Introduction**

Pathway enrichment analysis is an essential step for interpreting high-throughput (*omics*) data that uses current knowledge of genes and biological processes. A common application determines statistical enrichment of molecular pathways, biological processes and other functional annotations in long lists of candidate genes[1,2]. Genomic, transcriptomic, proteomic and epigenomic experiments emphasize distinct and complementary aspects of underlying biology and are best analyzed integratively, as is now routinely done in large-scale projects such as The Cancer Genome Atlas (TCGA)[3], Clinical Proteome Tumor Analysis Consortium (CPTAC), International Cancer Genome Consortium (ICGC)[4], Genotype-Tissue Expression (GTEx)[5] and others. Thus, simultaneous analysis of multiple candidate gene lists for characteristic pathways is increasingly needed. Numerous approaches are available for interpreting single gene lists. For example, the GSEA algorithm can detect up- and down-regulated pathways in gene expression datasets[6]. Web-based methods such as Panther[7], ToppCluster[8] and g:Profiler[9] detect significantly enriched pathways amongst ranked or unranked gene lists and are generally applicable to genes and proteins from various analyses. Some approaches allow analysis of multiple input gene lists however these primarily rely on visualization rather than data integration to evaluate the contribution of distinct gene lists towards each detected pathway[8,9]. Finally, no methods are available for unified pathway analysis of coding and non-coding mutations from whole-genome sequencing (WGS) data, or integrating these with other types of DNA aberrations such as copy number changes and balanced genomic rearrangements. We report the development of the ActivePathways method that uses data fusion techniques to address the challenge of integrative pathway analysis of multi-omics data. We demonstrate the method by analyzing known and candidate cancer driver genes with coding and non-coding somatic mutations in 2,583 whole cancer genomes of the ICGC-TCGA PCAWG project[10,11], prognostic pathways in breast cancer subtypes, and regulatory networks of tissue transcriptomes using the GTEx[5] compendium.

Characterization of genes and somatic mutations that drive oncogenesis is a central goal of cancer genomics research. Cancer genomes are characterized by few frequently mutated pan-cancer drivers such as *TP53*, less-frequent drivers with primarily tissue-specific effects and numerous infrequently mutated genes often referred to as *the long tail*. The majority of currently known driver mutations affect protein-coding sequence[12] and only few high-confidence non-coding drivers have been found, such as the mutation hotspots in the *TERT* promoter[13]. Discovery of non-coding driver mutations is a major goal of large cancer whole genome sequencing efforts such as PCAWG[10,11]. Pathway and network analysis of cancer mutations is a powerful approach

60   that uses knowledge of coding driver genes and their pathway annotations as priors to assist in

61   detection of weak driver variants including those in the non-coding genome[1]. The PCAWG project

62   has produced a consensus dataset of predicted protein-coding driver genes (CDS) and non-

63   coding regions of 5' and 3' untranslated elements (UTRs), promoters and enhancers of protein-

64   coding genes across 2,583 whole cancer genomes of multiple cancer types[14]. Driver gene p-

65   values in the dataset reflect the frequency and functional impact of somatic single nucleotide

66   variants (SNVs) and small insertions-deletions (indels) in these protein-coding and non-coding

67   genomic regions. Here we used our ActivePathways method to interpret these driver predictions

68   with pathway information including biological processes of Gene Ontology[15] and molecular

69   pathways defined by Reactome[16]. Two further case studies focused on prognostic molecular

70   pathways of breast cancer through integration of genomic and transcriptional alterations, and

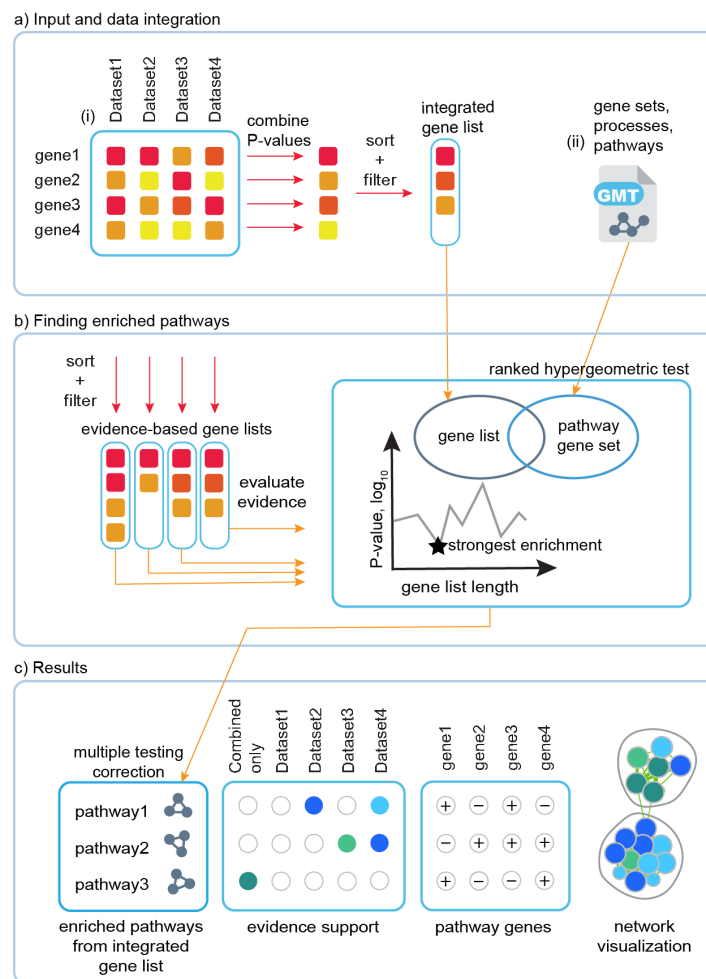71   gene regulatory networks associated with organ growth control in healthy human tissues.

72

73   **Results**

74   **Multi-omics pathway enrichment analysis with ActivePathways**

75   ActivePathways is a simple four-step method that extends our earlier work[9] (**Figure 1**). It requires

76   two input datasets. First, a table of *gene p-values* contains multiple p-values for every gene

77   representing different types of evidence such as gene significance in distinct omics experiments.

78   These could include p-values evaluating the significance of differential gene expression in tissues

79   of interest, gene essentiality, mutation or copy number alteration burden, and many others.

80   Second, a collection of *gene sets* represents molecular pathways, biological processes and other

81   gene annotations we refer to as *pathways*. Depending on the hypothesis, pathways may also

82   include other types of gene sets such as targets of transcription factors or microRNAs. In the first

83   step of ActivePathways, we derive an integrated gene list that aggregates significance from all

84   types of evidence for each input gene. The integrated gene list is compiled by fusion of gene

85   significance from different types of evidence using the Brown's extension[17] of the Fisher's

86   combined probability test, which conservatively adjusts for overall correlations of p-values in

87   estimating the overall significance of every gene. The integrated input gene list is then ranked by

88   decreasing significance and filtered using a lenient cut-off to capture a long tail of candidate genes

89   and to filter the bulk of insignificant ones (unadjusted $P_{gene}<0.1$). The integrated gene list is

90   analyzed with a ranked hypergeometric test for each pathway to capture smaller pathways tightly

91   associated with few top-ranking genes and broader processes with abundant albeit weaker

92   signals from larger subsets of input genes. The stringent family-wise multiple testing correction

93    method by Holm[18] is applied across pathways to reduce false positives ($Q_{pathway}$<0.05). In the third

94    step, candidate gene lists corresponding to distinct types of evidence are separately evaluated

95    using the above procedure. This step determines which pathways are significantly supported by

96    each of the input omics datasets and also reveals corresponding genes in each pathway.

97    Importantly, the step also highlights pathways that are only found through data integration and

98    are not apparent in any single type of omics evidence alone. In the fourth step, the method

99    provides input files for Enrichment Map[19] for visualizing and reducing the redundant set of all

100    detected pathways to a narrower, focused network of biological themes.

101



102    *Figure 1: Method overview. (a) ActivePathways requires as input (i) a matrix of gene P-values for different omics*

103    *datasets, and (ii) a collection of gene sets corresponding to biological pathways and processes. Gene p-values are*

104    *merged and filtered to produce an integrated gene list that combines evidence from omics datasets and is ranked by*

105    *decreasing significance with a lenient threshold. (b) Pathway enrichment analysis is conducted on the integrated gene*

106    *list as well as lists from individual omics datasets using the ranked hypergeometric test that determines the optimal*

107    *level of enrichment in the ranked gene sub-list for every pathway. (c) Pathways enriched in the integrated gene list are*

108    *corrected for multiple testing and significant findings are reported as results. Pathways enriched in individual omics*

109    *datasets are labelled by supporting evidence (colored nodes), and pathways only enriched in the integrated gene list*

110    *are highlighted separately. Pathway genes with significant signals in different omics data are also shown. Finally,*

111    *datasets of enriched pathways provided by ActivePathways are visualized as enrichment maps in Cytoscape where*

112    *nodes correspond to pathways and pathways with many shared genes are connected into networks representing*
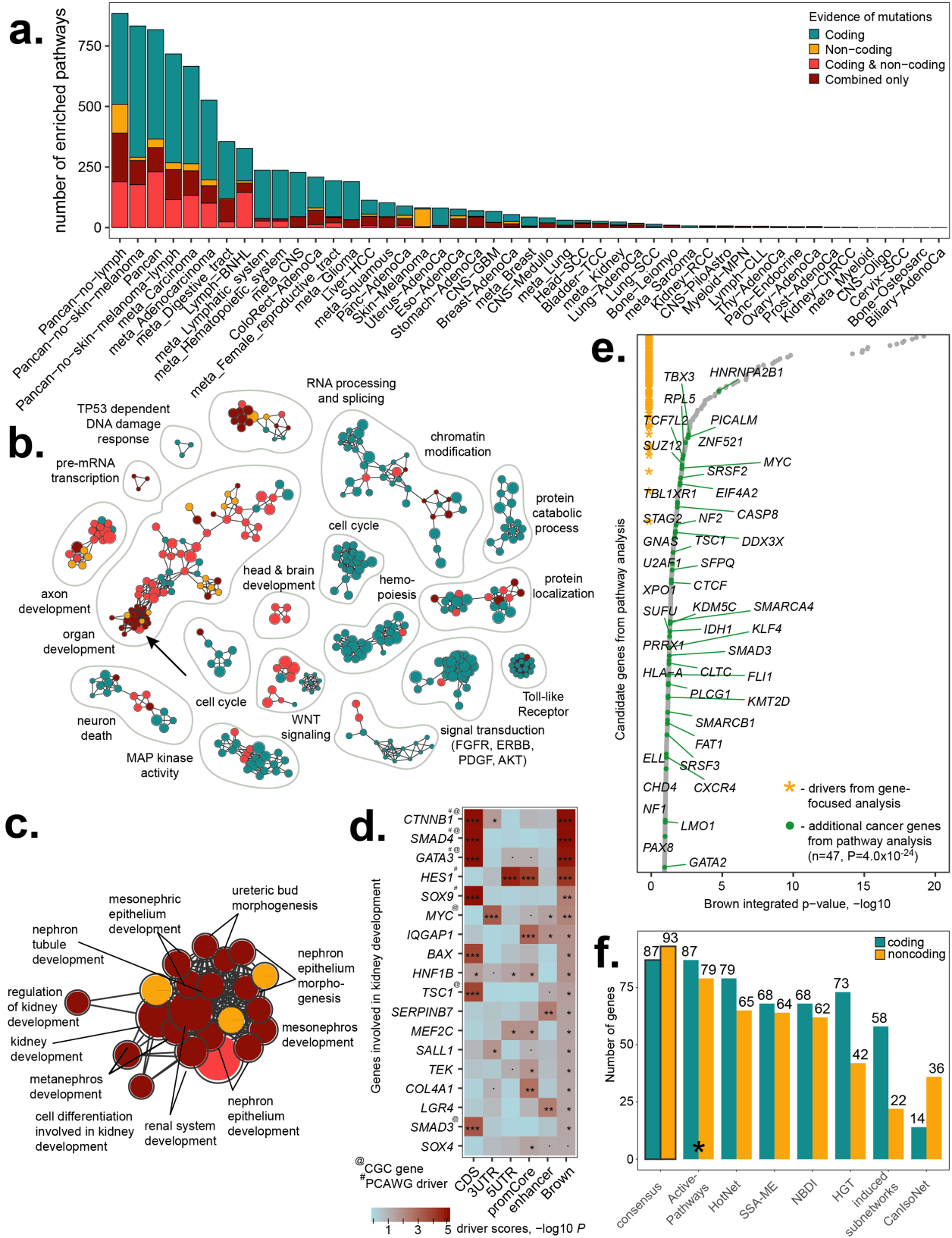
113    *broader biological themes.*

114

## Pathway analysis of coding and non-coding mutations in 2,500 whole cancer genomes

116    We performed integrative pathway analysis of coding and non-coding driver predictions across

117    29 cancer patient cohorts of histological tumor types and 18 meta-cohorts combining multiple

118    types of tumors, with 47 cohorts in total (**Supplementary Table 1**). ActivePathways found at least

119    one significantly enriched process or pathway in the majority of these cohorts (42/47 or 89%,

120    $Q_{pathway}$<0.05) (**Figure 2a**). We analyzed the omics evidence supporting predictions of enriched

121    pathways and found that most cohorts showed enrichments in pathways supported by protein-

122    coding driver scores of genes (37/47 or 79%). This serves as a positive control since the majority

123    of currently known cancer driver genes have frequent protein-coding mutations.

124    Non-coding mutations in genes also contributed to the discovery of frequently mutated biological

125    processes and pathways: 24/47 cohorts (51%) showed significantly enriched pathways that were

126    apparent when only analyzing non-coding driver scores separately for UTRs, promoters or

127    enhancers. The majority of cohorts (41/47 or 87%) revealed enriched pathways that were

128    apparent in the integrated gene list but not in any gene lists ranked by element-specific driver

129    scores, emphasizing the value of our integrative approach. As expected, cohorts with more patient

130    tumor samples generated more significantly enriched pathways (Spearman $\rho$=0.74, $P$=2.3x10$^{-9}$;

131    **Supplementary Figure 1**), suggesting that larger datasets are better powered to distinguish

132    rarely mutated genes involved in biological pathways and processes. Discovery of pathways

133    enriched in non-coding mutations suggests that pathway analysis is an attractive strategy for

134    illuminating the dark matter of the non-coding cancer genome.

ActivePathways: Paczkowska, Barenboim, *et al.*

136  *Figure 2. Pathway enrichment analysis of cancer driver genes with ActivePathways. (a) We analyzed consensus*
137  *driver genes with frequent somatic mutations by integrating mutation scores of protein-coding and non-coding*
138  *sequences (promoters, enhancers, and untranslated regions) across 47 cohorts of cancer patients with whole genome*
139  *sequencing data from tumors. Bar plot shows number of significantly enriched pathways (Q<0.05) stratified by*
140  *supporting evidence from driver predictions. The majority of pathways detected by ActivePathways are supported by*
141  *protein-coding mutations, as expected (dark green bars), while non-coding mutations (orange, red) reveal additional*
142  *pathways. Pathways shown in dark red are found only in the integrated gene list of coding and non-coding mutations*
143  *but not in gene lists of individual mutation scores. (b) Enrichment map shows groups of statistically significant pathways*
144  *characteristic of mutated genes in the adenocarcinoma cohort of 1,773 tumors. Nodes in the network diagram represent*
145  *pathways that are connected with edges if the pathways are similar and share many genes. Groups of similar pathways*
146  *were annotated manually. Nodes are colored by supporting evidence from coding and non-coding cancer mutations.*
147  *(c) The group of enriched kidney developmental processes is apparent from integrated evidence of coding and non-*
148  *coding mutations but is not found among coding or non-coding candidate genes separately (indicated with arrow in*
149  *enrichment map). (d) P-value heatmap shows driver scores of genes involved in kidney developmental processes*
150  *ranked by combined p-values of the integrated gene list (rightmost column). Top genes are expectedly detected as*
151  *significantly mutated driver genes in the PCAWG consensus list while additional pathway-derived genes of the long tail*
152  *of infrequent mutations are highlighted as well. Genes listed in the Cancer Gene Census (CGC) database are indicated*
153  *with @-symbol. (e) Integrated list of adenocarcinoma candidate driver genes used in the pathway enrichment analysis*
154  *includes the majority of driver genes detected in the gene-focused consensus analysis by PCAWG (orange asterisks)*
155  *and a long tail of infrequently mutated genes ranked by decreasing significance. Additional known cancer genes*
156  *detected in the pathway analysis are indicated with green dots and occur more frequently than expected from chance*
157  *alone. (f) Comparison of ActivePathways with six additional pathway and network analysis methods used in the*
158  *PCAWG project. ActivePathways best recovers the consensus lists of pathway-implicated driver (PID) genes with*
159  *coding and non-coding mutations. The consensus lists are shown in the leftmost bars of the plot and have been*
160  *compiled through a majority vote of the seven methods in the PCAWG pathway and network analysis working group.*

161  We studied the adenocarcinoma meta-cohort with 1,773 samples of 16 tumor types whose
162  integrated list of 432 candidate genes (unadjusted $P_{gene}$<0.1) associated with 526 significantly
163  enriched pathways ($Q_{pathway}$<0.05) (**Figure 2b**). As expected, the majority of pathways were only
164  supported by genes with frequent coding mutations (328/526 or 62%). However, 101 pathways
165  were supported by both coding and non-coding gene mutations, 72 were only apparent in the
166  integrated analysis of all evidence, and 25 were only found among genes with significant non-
167  coding mutations, thus expanding the set of candidate driver mutations in the non-coding cancer
168  genome and demonstrating the value of integrative pathway analysis.

169  The major biological themes with frequent protein-coding mutations included hallmark cancer
170  processes like *apoptotic signaling pathway* (24 genes; $Q_{pathway}$=4.3x10$^{-5}$) and *mitotic cell cycle* (8
171  genes; $Q_{pathway}$=0.0026), and additional biological processes such as chromatin modification and
172  RNA splicing that are increasingly recognized in cancer biology. Thus, our method captures the
173  expected cancer pathways among driver genes with protein-coding mutations as positive controls.

174    In contrast to these solely protein-coding driver associations, a large group of developmental

175    processes and signal transduction pathways was enriched in genes with coding as well as non-

176    coding mutations; for example *embryo development process* was supported by mutations in

177    exons, 3'UTRs and gene promoters (68 genes; $Q_{pathway}$=2.9x10$^{-12}$), while *repression of WNT target*

178    *genes* was only apparent in the integrated analysis of coding and non-coding mutations but not

179    in either alone (5 genes, $Q_{pathway}$=0.016; REAC:4641265). Thus, our method evaluates

180    contribution of omics evidence towards pathway enrichments and finds additional associations

181    that are not apparent in any provided dataset.

182

**ActivePathways highlights pathway-associated cancer genes in the long tail of infrequent**

**non-coding mutations**

185    We focused on a group of processes involved in kidney development that were only detected in

186    the integrated analysis (**Figure 2c-d**). ActivePathways found 18 genes involved in these

187    processes, only five of which were predicted as driver genes in the consensus driver analysis of

188    the PCAWG project[14]. Additional known cancer genes included the oncogene *MYC* with 13

189    patients with 3'UTR mutations ($P_{UTR3}$=4.8x10$^{-4}$; $Q_{UTR3}$=0.42), the transcription factor *SMAD3* of

190    the TGF-β pathway with 14 patients with protein-coding mutations ($P_{CDS}$=4.0x10$^{-4}$; $Q_{CDS}$=0.37)

191    and the growth inhibitory tumor suppressor gene *TSC1* with 23 patients with protein-coding

192    mutations ($P_{CDS}$=1.4x10$^{-4}$; $Q_{CDS}$=0.17) as well as candidate cancer genes such as *IQGAP1* with

193    10 patients with promoter mutations ($P_{promoter}$=8.2x10$^{-4}$; $Q_{promoter}$=0.62) that encodes a signaling

194    protein that regulates cell motility and morphology. The additional genes remained below the

195    FDR-adjusted significance cut-off in the gene-focused consensus driver analysis, however were

196    found by ActivePathways due to pathway associations with frequently mutated developmental

197    genes. These results highlight the potential of our method to find known and candidate cancer

198    genes with rare coding and non-coding driver mutations through pathway-driven data integration.

199    We evaluated 333 candidate driver genes from the pathway analysis of the adenocarcinoma

200    cohort (**Figure 2e**). These included as positive controls 60/64 significantly mutated genes

201    identified in the PCAWG consensus driver analysis[14], and an additional 47 genes of the COSMIC

202    Cancer Gene Census database[12], significantly more than expected by chance alone (seven

203    genes expected, Fisher's exact $P$=4.0x10$^{-24}$), including *MYC, IDH1, NF1,* and *BCL9*. Additional

204    genes were detected for several reasons. First, the integrated gene list was filtered using a lenient

205    statistical cut-off ($P_{gene}$<0.1) compared to a more stringent gene-focused driver analysis

206    ($Q_{gene}$<0.05). This resulted in 273/333 pathway-associated genes of the long tail that remained

207 below the significance threshold in the driver analysis. Second, the integration procedure
208 combined multiple weaker p-values (coding regions, promoters, UTRs, enhancers) to a single
209 stronger p-value for 17/333 pathway-associated genes including six cancer genes (*HNRNPA2B1*,
210 *STAG2*, *TCF7L2*, *SUZ12*, *CLTC*, *ZNF521*) and improved the overall ranking of 220/333 genes
211 among the input data, better explaining their membership in pathways and processes. However,
212 a majority of all genes showed reduced significance after the integration procedure and were
213 excluded from the pathway analysis, as the Brown combined p-value remained below the
214 significance cut-off compared to any individual p-values of mutations in coding and non-coding
215 regions of genes (3,112/3,543 or 88% genes with unadjusted $min(P_{gene})$<0.1 showed unadjusted
216 Brown $P_{gene}$>0.1). Fourth, the evidence evaluation step of the method identified pathway
217 enrichments in gene lists ranked by individual sources of evidence and highlighted additional
218 genes that did not pass significance cut-offs of the integration procedure. Thus, ActivePathways
219 finds additional cancer genes in the long tail of mutations that are highlighted due to their pathway
220 associations but remain below the significance cut-off in the gene-by-gene analysis.

221

## Benchmarking demonstrates the robustness and sensitivity of ActivePathways

223 We carefully benchmarked ActivePathways using multiple approaches. First, we compared its
224 performance with six diverse methods used in the PCAWG pathway and network analysis working
225 group[20] (Hierarchical HotNet[21,22], SSA−ME[23], NBDI[24], induced subnetwork analysis[22],
226 CanIsoNet[Kahraman et al, in prep], and hypergeometric test). The methods used molecular pathway and
227 network information to analyze the PCAWG dataset of predicted cancer driver genes[14], and a
228 subsequent consensus procedure derived pathway-implicated driver (PID) gene lists with coding
229 (PID-C) and non-coding (PID-N) mutations based on a majority vote. Our method recovered PID-
230 C and PID-N gene lists with the highest accuracy: 100% of coding driver genes (87/87) and 85%
231 of non-coding candidates (79/93) were detected (**Figure 2f**).

232 We evaluated the robustness of ActivePathways to parameter variations and missing data. We
233 varied the parameter $P_{gene}$ that determines the ranked gene lists used in the pathway enrichment
234 analysis (default threshold $P_{gene}$<0.1). The majority of cohorts (40/47 or 85%) retrieved
235 significantly enriched pathways even with a considerably more stringent threshold ($P_{gene}$<0.001),
236 however 67% fewer pathways were found compared to the default threshold in the median cohort
237 (**Supplementary Figure 2**). We then evaluated the robustness of ActivePathways to missing data
238 by randomly removing subsets of driver scores from the initial dataset. Even when removing 50%
239 of gene driver scores with $P$<0.001, the majority of cohorts (37/47 or 79%) were found to have at

240   least one significantly enriched pathway however 66% fewer pathways were found on average

241   (**Supplementary Figure 3**).

242   We tested ActivePathways with data simulations through 1,000 datasets for each of 47 patient

243   cohorts and found no significant pathways in 92% of simulations (**Supplementary Figure 4**).

244   Simulated data were obtained by randomly reassigning driver scores to different genomic

245   elements, a conservative approach that disrupts gene and pathway annotations while retaining

246   strong scores in the data. The median family-wise false discovery rate across cohorts (7.2%)

247   slightly exceeded the applied multiple testing correction ($Q<0.05$). Higher rates were observed in

248   cohorts including melanoma tumors, potentially due to abundant promoter mutations caused by

249   impaired nucleotide excision repair in protein-bound genomic regions[25]. We evaluated quantile-

250   quantile (QQ) plots of pathway-based p-values from ActivePathways and found that p-values from

251   observed gene scores often deviated from the expected uniform distribution and appeared

252   statistically inflated (**Supplementary Figure 5**). However, p-values derived from simulated gene

253   scores showed no inflation in our simulations. Anticipating that the strongest cancer driver scores

254   associate with protein-coding sequence, we studied datasets with simulated protein-coding gene

255   scores and true non-coding scores. As expected, these partially simulated datasets expectedly

256   showed less p-value inflation, suggesting that highly significant known cancer genes involved in

257   many different pathways are responsible for inflation. Statistical testing of highly redundant

258   pathways and processes violates the independence assumption of statistical tests and multiple

259   testing procedures, a known caveat of pathway enrichment analysis[1,2], which likely explains the

260   observed distribution of significance values of our method.

261   Collectively, these benchmarks show that ActivePathways is a sensitive and robust method for

262   detecting significantly enriched pathways and processes through integrative analysis of

263   multivariate omics data.

264

265   **Clinical analysis of genomic and transcriptional alterations of breast cancer subtypes**

266   **reveals prognostic value of apoptotic, immune response and ribosomal genes**

267   To demonstrate an integrative analysis of patient clinical information with multiple types of omics

268   data , we then studied the pathways and processes associated with patient prognosis in breast

269   cancer. We leveraged the METABRIC dataset[26] using 1,780 breast cancer samples drawn from

270   all four subtypes (HER2-enriched, basal-like, luminal-A, luminal-B) and evaluated all genes using

271   three types of prognostic evidence. Gene expression profiles were deconvolved as mRNA

272  abundance levels in tumor cells (TC) and tumor-adjacent cells (TAC) using the ISOpure

273  algorithm[27] and associated with these data with patient survival using median dichotomization and

274  log-rank tests. Gene copy number alterations (CNA) were included as the third type of evidence

275  and associated with patient survival using log-rank tests.

276  ActivePathways highlighted 192 significantly enriched GO biological processes and Reactome

277  pathways across the four subtypes, of which nine were enriched in multiple subtypes and 33 were

278  only apparent through the integrative pathway analysis but not in any omics evidence alone.

279  Enrichment maps of significant results revealed immune response, apoptosis, ribosome

280  biogenesis and chromosome segregation as the major groups of prognosis-associated pathways

281  (**Figure 3a**).

282  Immune activity was associated with prognostic genes in basal-like and HER2-enriched breast

283  cancers with significant enrichment of GO processes such as immune system development

284  ($Q_{basal}$=3.0x10$^{-4}$, 113 genes; $Q_{HER2}$=0.035, 61 genes) and lymphocyte differentiation

285  ($Q_{HER2}$=6.8x10$^{-4}$, 46 genes; $Q_{basal}$=8.4x10$^{-4}$, 45 genes). The majority of genes of immune system

286  development were associated with improved patient prognosis upon increased gene expression

287  in tumor cells or tumor-adjacent cells, comprising 50/61 genes in the HER2-enriched subtype and

288  78/113 genes in the basal subtype (**Figure 3b**). Interestingly, only a minority of these genes (10)

289  were significant in both of the two subtypes, suggesting different modes of immune activity in

290  subtypes and emphasizing the power of our pathway-based approach. Basal-like breast cancers

291  were associated with additional 67 terms involving immune response and blood cells, however

292  no immune related terms were enriched for luminal subtypes of breast cancers. Prognostic

293  features of immune-related genes in HER2-enriched and basal-like breast cancers are well

294  known[28,29]. Our pathway-based findings indicate that immune activity in breast tumor cells and in

295  the surrounding microenvironment negatively affect tumor progression and benefits the patient.
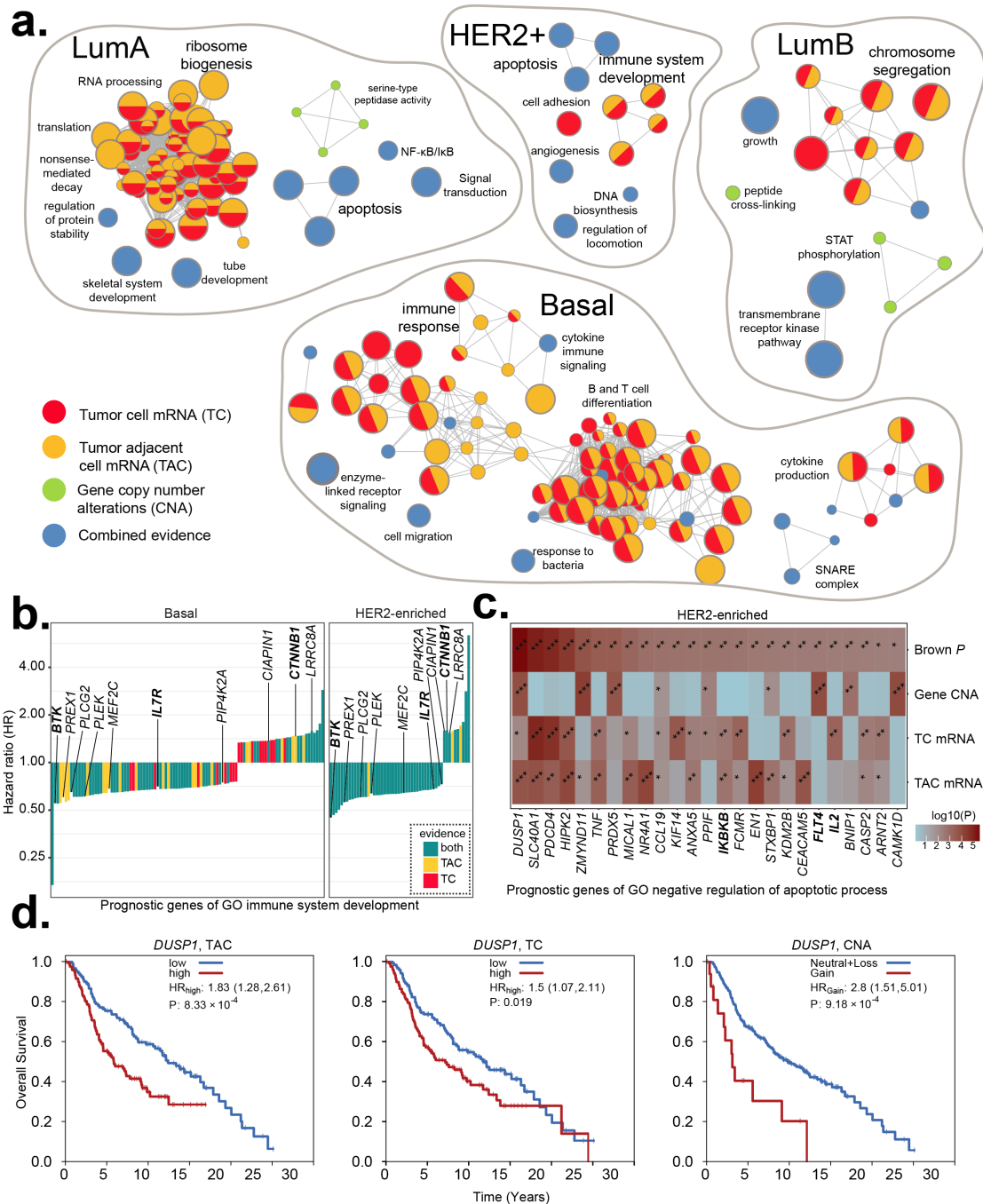
296  Apoptosis was associated with patient prognosis in HER2-enriched and luminal-A breast cancers

297  through enriched GO processes such as negative regulation of apoptotic process ($Q_{HER2}$=0.030,

298  122 genes; $Q_{luminalA}$=0.015, 228 genes) and programmed cell death ($Q_{HER2}$= 0.015, 125 genes;

299  $Q_{luminalA}$= 0.016, 231 genes) (**Figure 3c**). Anti-apoptotic pathways were only detected in the

300  integrative analysis and not in genomic and transcriptomic gene signatures separately. Among

301  the genes negatively regulating apoptosis, *DUSP1* provided the strongest prognostic signal in

302  HER2-enriched breast cancers. This was apparent in the molecular stratification of samples by

303  mRNA of tumor cells (log-rank $P_{TC}$=0.019, HR=1.5) and tumor-adjacent cells ($P_{TAC}$=8.3x10$^{-4}$,

304  HR=1.83) as well as gene copy number amplifications ($P_{CNA}$=9.8x10$^{-4}$, HR=2.8) (**Figure 3d**).

305  *DUSP1* encodes a phosphatase signaling protein of the MAPK pathway that is over-expressed in

306  malignant breast cancer cells and inhibits apoptotic signaling[31]. *HER2* over-expression is known

307  to suppress apoptosis in breast cancer[30]. Anti-apoptotic signaling is a hallmark of cancer and

308  expectedly associated with worse patient prognosis.

309  ActivePathways also identified prognostic pathway associations in single subtypes of breast

310  cancer. For example, the prognostic genes for luminal-B subtype were enriched for chromosome

311  segregation ($Q_{luminalB}$=0.017, 41 genes) and related biological processes of GO. In agreement with

312  this finding, problems with chromosome segregation have been associated with worse outcome

313  in breast cancer[32]. As another example, luminal-A breast cancers were associated with prognosis

314  in ribosomal and RNA processing genes, such as ribosome biogenesis ($Q_{luminalA}$=6.9x10$^{-10}$, 60

315  genes), and rRNA metabolic process ($Q_{luminalA}$=1.8x10$^{-13}$, 64 genes). Although not described

316  specifically in the luminal-A subtype, ribosomal mRNA abundance has been shown to be

317  prognostic in breast cancer as a marker of cell proliferation[33,34]. In summary, ActivePathways can

318  be used for integrating clinical data with multi-omics information of molecular alterations. Such

319  analyses can provide leads for functional studies and biomarker development.

**Figure 3. Prognosis-associated pathways in four molecular subtypes of breast cancer.** *(a) Enrichment maps of prognostic pathways and processes were found in an integrative analysis of mRNA abundance in tumor cells (TC), tumor-adjacent cells (TAC) and gene copy number alterations (CNA). Multi-colored nodes indicate pathways that were prognostic according to several types of molecular evidence. Blue nodes indicate pathways that were only apparent through merging of molecular signals. (b) Hazard ratios (HR) of prognostic genes of immune system development in basal and HER2-enriched subtypes of breast cancer. Strongest HR of TC, TAC is shown. Genes commonly found in basal and HER2-enriched tumors are shown. (c) Heatmap shows genes and corresponding p-values of the GO process "negative regulation of apoptotic process" found as prognostic in HER2-enriched breast cancer. Top row of the heatmap shows Brown p-values of merged evidence. (d) Kaplan-Meier plots show the strongest prognostic signal of the above apoptotic process associated with the DUSP1 encoding a protein phosphatase. DUSP1 significantly associates with reduced patient survival through increased tumor-adjacent mRNA level (left), increased tumor mRNA level (center) and gene copy number amplification (right). Known cancer genes are shown in boldface letters.*
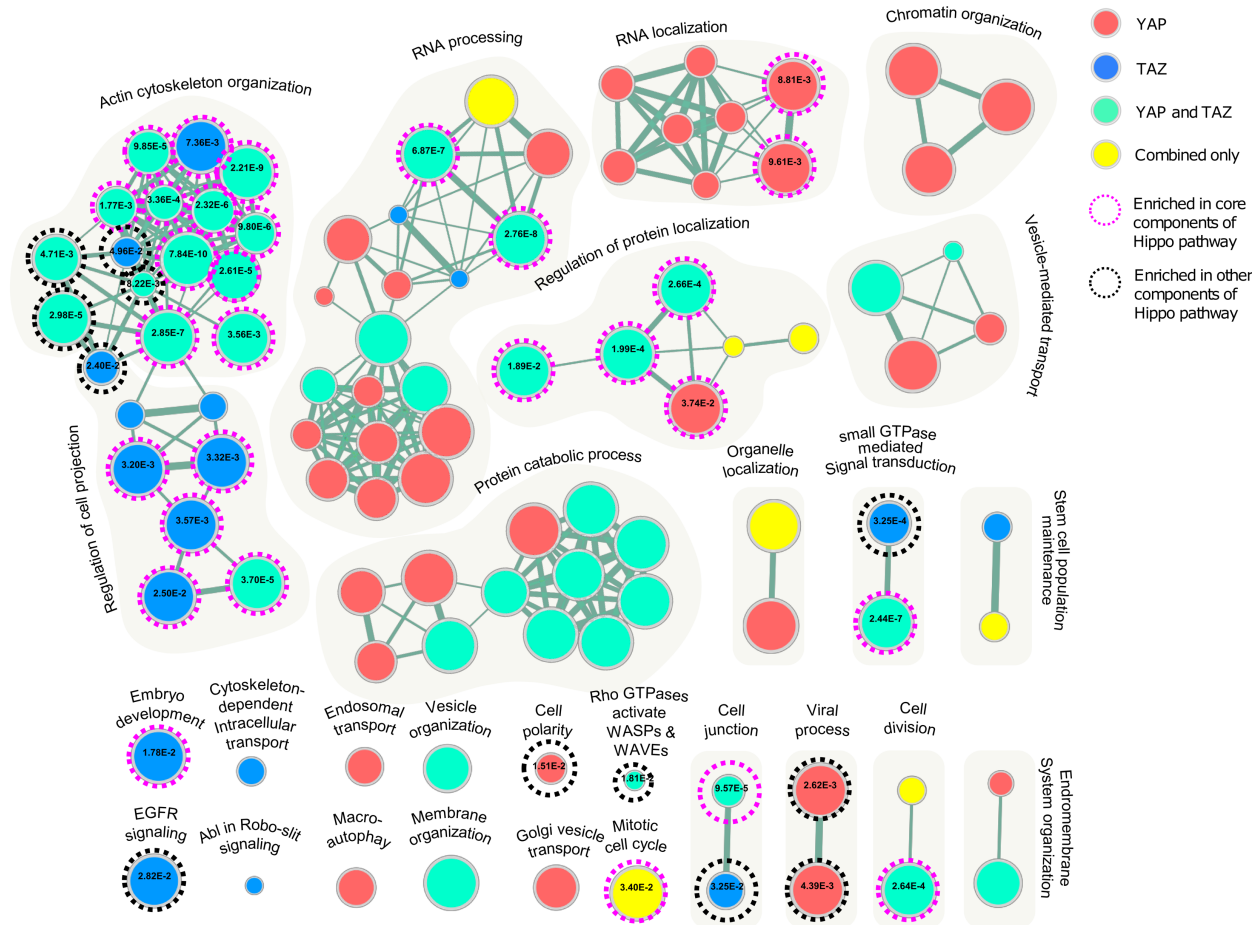
**333  Co-expression analysis of Hippo master regulators across 54 human tissues recovers**

**334  associated biological processes and genes**

335  To study the use of ActivePathways in the context of healthy human tissues, we analyzed the

336  dataset of 11,688 transcriptomes of 54 tissues from the GTEx project[5], focusing on the Hippo

337  signaling pathway involved in organ size control, tissue homeostasis and cancer[35,36]. We studied

338  gene co-expression networks downstream of YAP and TAZ, the two master transcription factors

339  of Hippo signaling, encoded by *YAP1* and *WWTR1*. YAP and TAZ are the evolutionarily

340  conserved key effectors of the Hippo signaling in mammals. Inhibition of YAP/TAZ-mediated

341  transcription regulates organ size control and tissue homeostasis in response to a wide range of

342  intracellular and extracellular signals including cell-cell interactions, cell polarity, mechanical cues,

343  ligands of G-protein-coupled receptors, and cellular energy status. We retrieved 2,117 putative

344  Hippo transcriptional target genes that showed significant positive co-expression with either or

345  both of the transcripts of *YAP* and *TAZ* across the human tissues in the GTEx dataset

346  ($Q_{gene}$<0.05). We used a robust rank aggregation method[37] and retrieved transcriptional targets

347  that were co-expressed with YAP or TAZ in a relatively large number of human tissues.

348  Analysis of the target genes using ActivePathways resulted in 101 significantly enriched pathways

349  ($Q_{pathway}$<0.05), including 39 supported by both sets of target genes, 37 supported by YAP1

350  targets, 18 supported by TAZ targets, and seven only apparent in the integrated list of target

351  genes (**Figure 4**). The major biological themes of pathways and processes included regulation of

352  cell polarity and cell junction, embryonic development, EGFR signaling, maintenance of stem cell

353  population, actin cytoskeleton, and rho GTPase signaling that are all directly or indirectly related

354  to Hippo signaling. We validated our analysis using 207 Hippo-related genes from review

355  papers[35,36] and confirmed that 83/101 pathways found by ActivePathways contained at least one

356  of 59 Hippo-related genes, while 41 pathways were significantly enriched in Hippo-related genes

357  ($Q$<0.05). However, the majority of genes documented in the literature (148/207) were not

358  detected in the pathway analysis, potentially due to their post-transcriptional regulation or tissue-

359  specific roles. Our analysis highlights known and candidate genes and pathways related to Hippo

360  signaling and showcases the use of ActivePathways for functional analysis of transcription

361  regulatory networks.

ActivePathways: Paczkowska, Barenboim, *et al.*



362

**Figure 4. Pathway enrichment analysis of Hippo co-expression targets across human tissues.** *Enrichment map of pathways characteristic of genes co-expressed with transcription factors YAP and TAZ of the Hippo pathway across human tissues in the GTEx dataset. The Hippo pathway is involved in organ growth control and its predicted target genes are enriched in related biological processes and pathways. Nodes represent significantly enriched pathways that are colored by supporting evidence from co-expression targets of YAP or TAZ (red, blue), both transcription factors (green) or only the integrated list of target genes (yellow). We validated the detected pathways using a list of Hippo-related genes compiled from recent review papers and found that the majority of detected pathways included Hippo-related genes and 40% of pathways were enriched in these genes (indicated with dotted circles, enrichment p-values shown in nodes).*

## Discussion

Integrative pathway enrichment analysis helps distill thousands of high-throughput measurements to a smaller number of pathways and biological themes that are most characteristic of the experimental data, ideally leading to mechanistic insights and novel candidate genes for follow-up studies. The primary advantage of our method is the fusion of gene significance across multiple

378    omics datasets. This allows us to identify additional pathways and processes that are not apparent

379    individually in any analyzed dataset. In our example of cancer driver discovery, pathway analysis

380    is complementary to gene-focused driver discovery as it also focuses on sub-significant genes

381    with coding and non-coding mutations clustered into known and novel biological processes of

382    cancer. In the clinical analysis of breast cancer subtypes, we find prognostic genes and pathways

383    active in tumor cells, the microenvironment, or both. A subset of these findings, such as anti-

384    apoptotic signaling, is only apparent through data integration.

385    Our general pathway analysis strategy is applicable to diverse kinds of omics datasets where

386    well-calibrated p-values are available for the entire set of genes or proteins. One may study a

387    series of genomic, transcriptomic, or proteomic experiments or combine these into a multi-omics

388    analysis. Data from epigenomic experiments and genome-wide association studies can be

389    analyzed after genome-wide signals have been appropriately mapped to genes. Clinical and

390    phenotypic information of patients can be also included through association and survival statistics.

391    Our method is expected to work with unadjusted as well as multiple-testing adjusted p-values,

392    however it is primarily intended for un-adjusted p-values for increased sensitivity. P-value

393    adjustment for multiple testing is conducted at the pathway level rather than at a gene level. P-

394    values from omics datasets are easier to interpret than raw signals as gene-based p-values are

395    expected to account for experimental and computational biases specific to each analyzed dataset,

396    while accounting for multi-omics factors comprehensively in a single generally applicable

397    pathway-based model would be likely impossible. In our example of cancer driver discovery,

398    appropriately computed p-values account for confounding factors of somatic mutations such as

399    gene sequence length and nucleotide content, mutation signatures active in different types of

400    tumors[38] and biological cofactors of mutation frequency such as transcription and replication

401    timing[39], while pathway analysis of mutation counts or frequencies would maintain such biases in

402    results.

403    Our analysis comes with important caveats. First, we only evaluate genes annotated in pathway

404    databases that have variable coverage, rely on frequent data updates[40] and may miss novel

405    sparsely annotated candidate genes. The most general pathway enrichment analysis considers

406    biological processes and molecular pathways however many kinds of gene sets available in

407    resources such as MSigDB[41] can be used to expand the scope of ActivePathways. Second,

408    pathway information is highly redundant and analysis of rich *omics* datasets often results in many

409    significant results reflecting the same underlying pathway. We address this redundancy by

410    visualizing and summarizing pathway results as enrichment maps[2,19] that help distill general

411 biological themes comprised of multiple similar pathways and processes. Statistical inflation of
412 results accompanied by biological redundancy is addressed by a stringent multiple testing
413 correction. Third, the analysis treats pathways as gene sets and does not consider their
414 interactions. This expands the scope of our analysis to a wider repertoire of pathways and
415 processes as reliable mechanistic interactions are often context-specific and limited to a small
416 subset of well-studied signaling pathways. Several methods such as HotNet[21], PARADIGM[42] and
417 GeneMania[43] model pathways and *omics* datasets through gene and protein interactions.

418 Translation of discoveries into improved human health through actionable mechanistic insights,
419 biomarkers, and molecular therapies is a long-standing goal of biomedical research. Next-
420 generation projects such as ICGC-ARGO (https://www.icgcargo.org/) aim to collect multi-*omics*
421 datasets with detailed clinical profiles of patients and thus present novel challenges for pathway
422 and network analysis techniques. In summary, ActivePathways is integrative pathway analysis
423 method that improves systems-level understanding of cellular organization in health and disease.

424 **Methods**

425 **Integrated and evidence-based gene lists.** The main input of ActivePathways is a matrix of p-
426 values where rows include all genes of a genome and columns correspond to omics datasets. To
427 interpret multiple omics datasets, a combined p-value was computed for each gene using a data
428 fusion approach, resulting in an integrated gene list. The integrated gene list was computed gene-
429 by-gene by merging all p-values of a given gene into one combined p-value using the Brown's
430 extension[17] of the Fisher's combined probability test that accounts for overall co-variations of p-
431 values from different sources of evidence. The integrated gene list of Brown p-values was ranked
432 in order of decreasing significance and filtered using a lenient threshold of unadjusted $P<0.1$.
433 Evidence-based gene lists representing different omics datasets were based on ranked P-values
434 from individual columns of the input matrix, using the same significance threshold.

435 **Statistical enrichment of pathways.** Statistical enrichment of pathways in significance-ranked
436 lists of candidate genes was carried out with the ranked hypergeometric test. The test considered
437 one pathway gene set at a time and analyzed increasing subsets of input genes from the top of
438 the ranked gene list. The same procedure was used for integrated and evidence-based gene lists.
439 At each iteration, the test computed the hypergeometric enrichment statistic and P-value for the
440 set of genes shared by the pathway and top sub-list of the input gene list. For optimal processing
441 speed, only gene lists ending with a pathway-related gene were considered as these most impact
442 significance of enrichment. The ranked hypergeometric statistic selected the input gene sub-list

443 that achieved the strongest enrichment and the smallest p-value as the final result for the given

444 pathway, as

445
$$(P_{pathway}, G) = \{\min, \underset{n}{\arg\min}\} \sum_{x=k}^{min(n,K)} \frac{\binom{K}{k}\binom{N-K}{n-k}}{\binom{N}{n}}$$

446 where $P_{pathway}$ stands for the hypergeometric P-value of the pathway enrichment at the optimal

447 sub-list of the significance-ranked candidate genes, $G$ represents the length of the optimal sub-

448 list, i.e. the number of top genes from the input gene list, $N$ is the number of protein-coding genes

449 with annotations in the pathway database, i.e., in Gene Ontology and Reactome, $K$ is the total

450 number of genes in a given pathway, $n$ is the number of genes in a given gene sub-list considered,

451 and $k$ is the number of pathway genes in the considered sub-list. For a conservative estimate of

452 pathway enrichment, we considered as background $N$ the universe of genes contained in pathway

453 databases and ontologies rather than the complete repertoire of protein-coding genes. To obtain

454 candidate genes involved in the pathway of interest, we intersected pathway genes with the

455 optimal sub-list of candidate genes. The ranked hypergeometric p-value was computed for all

456 pathways and resulting p-values were corrected for multiple testing using the conservative Holm-

457 Bonferroni family-wise error rate (FWER) method[18]. Significant pathways were reported ($Q<0.05$).

458 **Evaluating *omics* evidence of enriched pathways.** The integrated gene list was analyzed the

459 using ranked hypergeometric test and enriched pathways were reported as results. Each

460 evidence-based gene list representing an omics dataset was also analyzed for enriched pathways

461 with the ranked hypergeometric test. Pathways found in the integrated gene list were labelled for

462 supporting evidence if they were also found as significant in any evidence-based gene list. A

463 pathway was considered to be found only through data integration and labelled as *combined-only*

464 if it was identified as enriched in the integrated gene list but was not identified as enriched in any

465 of the evidence-based gene lists at equivalent significance cutoffs ($Q<0.05$). Each detected

466 pathway was additionally annotated with pathway genes apparent in the optimal sub-list of

467 candidate genes, separately for the integrated gene list and each evidence-based gene list.

468 *Gene scores of cancer mutations.* We analyzed p-values of genes reflecting their statistical

469 significance as candidate cancer drivers for multiple cohorts of cancer patients with whole

470 genome sequencing data. The scores were compiled in the driver discovery analysis of the

471 PCAWG project as a consensus of multiple independent methods[14]. The input matrix of gene

472 scores (*P*-values) included all protein-coding genes as rows and their genomic elements as

473 columns (exons, 5' and 3' untranslated regions (UTRs), promoters, enhancers). Elements with

474 missing p-values were assigned $P=1$. Genes with multiple enhancers were assigned the score of

475  the most significant enhancer, and enhancers with more than five associated genes were
476  excluded prior to selection.

***Pathways and processes.*** We used gene sets corresponding to biological processes of Gene Ontology[15] and molecular pathways of the Reactome database[16] downloaded from the g:Profiler web server[9]. Large general gene sets with more than a thousand genes and small specific gene sets with less than five genes were excluded.

**Enrichment map visualization.** ActivePathways provides input files for the EnrichmentMap app[19] of Cytoscape[45] for network visualization of similar pathways and their coloring according to supporting omics evidence. Enrichment maps for adenocarcinoma driver mutations, breast cancer prognostics, and Hippo transcriptional networks were visualized with stringent pathway similarity scores (Jaccard and overlap combined coefficient 0.6) and manually curated for the most representative groups of similar pathways and processes. Singleton pathways that were redundant with larger groups of pathways were discarded. Coloring of pathways in the adenocarcinoma enrichment map was rearranged by merging colors of pathways supported by non-coding mutation scores of promoters, enhancers and/or UTRs into one group.

**Analysis of coding and non-coding mutations of the PCAWG pan-cancer dataset.** We used ActivePathways to analyze driver predictions of coding and non-coding mutations across >2,500 whole cancer genomes of the ICGC-TCGA PCAWG Project. P-values of driver predictions were computed separately for protein-coding sequences, promoters, enhancers and untranslated regions (UTR3, UTR5) in the PCAWG driver discovery study by Rheinbay *et al*[14] across multiple subsets of samples representing histological tumor types and pan-cancer cohorts. We used gene-enhancer mapping predictions provided by PCAWG, excluded enhancers with more than five target genes, and selected the most significant enhancer for each gene, if any. Unadjusted p-values for coding sequences, promoters, enhancers and UTRs were compiled as input matrices and analyzed as described above. Missing p-values were interpreted as ones. Results from ActivePathways were validated with two lists of cancer genes. Predicted drivers from the gene-focused PCAWG driver analysis[14] were selected as statistically significant findings ($Q<0.05$) following a stringent multiple testing correction spanning all types of elements (exons, UTRs, promoter, enhancers). The curated list of known cancer genes was retrieved from the COSMIC Cancer Gene Census (CGC) database[12]. One-tailed Fisher's exact tests were used to estimate enrichment of these genes using all protein-coding genes as background.

**Analysis of prognostic genes in breast cancer.** ActivePathways was used to evaluate prognostic pathways in breast cancer using multiple types of omics data. mRNA gene expression

508    data and gene copy number alteration (CNA) data of the were derived from the METABRIC cohort

509    of 1,991 patients with a single primary fresh frozen breast cancer specimen each[26]. Curtis *et al*[26]

510    classified the patients into the intrinsic breast cancer subtypes using the PAM50 mRNA-based

511    classifier[44] resulting in 330 basal-like breast cancers, 238 HER2-enriched breast cancers, 721

512    luminal-A breast cancers, 491 luminal-B breast cancers. Using these data, we computationally

513    deconvolved tumor cell (TC) mRNA and tumor adjacent cell (TAC) mRNA abundance levels from

514    the bulk profiled specimens. TC mRNA was deconvolved using ISOpure[27] run on MATLAB

515    release 2010b. TAC mRNA was computed using the ISOpure.calculate.tac function from the R

516    package ISOpureR v1.1.2. ISOpure was run independently for each breast cancer subtype. The

517    mRNA univariate survival analysis was conducted as follows. For each gene, patients were

518    dichotomized based on mRNA abundance. Dichotomization was either based on the median

519    mRNA abundance for that gene or a fixed value of 6.5. Based on the mRNA abundance

520    distribution of genes on the Y chromosome in female samples, 6.5 was estimated as the threshold

521    for noise for non-expressed genes. Median dichotomization was used if the median was above

522    6.5 or if there were no events in one of the groups when dichotomizing based on 6.5. The high

523    and low mRNA abundance groups were compared by univariate log-rank tests for overall survival.

524    TC and TAC mRNA abundance were evaluated independently. Survival modelling was performed

525    in the R statistical environment (v3.4.3) using the survival package (v2.42-3). The CNA univariate

526    survival analysis was conducted as follows. For each gene, we assessed whether more gains or

527    losses were apparent. The copy number status with a higher count was subsequently used to

528    separate patients into two groups: those with the chosen copy number status and the remaining

529    patients. The two groups were then used for overall survival modelling with log-rank tests in the

530    R statistical environment (v3.4.3) using the survival package (v2.42-3).

531    **Co-expression analysis of GTEx transcriptomes.** The RNAseq dataset of human tissues was

532    downloaded from GTEx v7 data portal (https://www.gtexportal.org/home/). The dataset included

533    transcript abundance values of 21,518 protein-coding genes in 11,688 samples across 54 tissues.

534    Tissues with less than 25 available samples and low gene expression (mean TPM<1.0) were

535    excluded from further analysis. Positive pairwise Pearson correlations of gene expression values

536    of *YAP* and *TAZ* (symbols *YAP1, WWTR1*) and their putative target genes were investigated in

537    individual tissues and ranked by statistical significance of correlation tests. Tissue-specific ranked

538    correlations of target genes were then integrated into two master lists of target genes of YAP and

539    TAZ, respectively, reflecting target genes that were consistently positively co-regulated with

540    corresponding transcripts across a significant subset of considered human tissues. We used the

541    robust rank aggregation (RRA) method developed by Kolde *et al*[37] and filtered co-expressed

542 genes by significance using the default parameters of RRA ($Q_{gene}$<0.05). Significantly enriched

543 pathways among the putative target genes of YAP and TAZ were detected using ActivePathways.

544 We validated the pathways by investigating their agreement with known Hippo-related genes from

545 recent review papers[35,36]. We tested each pathway for enrichment of literature-derived Hippo

546 genes using Fisher's exact tests and filtered significant findings after multiple testing correction

547 (Q<0.05).

548 **Method benchmarking.** We benchmarked ActivePathways using multiple approaches, including

549 simulated datasets, parameter variations, and partial replacement of strong scores with missing

550 values. Benchmarking was carried out with the PCAWG dataset of coding and non-coding cancer

551 driver predictions. To evaluate false discovery rates of ActivePathways, we created simulated

552 datasets by randomly reassigning all observed driver scores to random genes and genomic

553 elements. Simulations were conducted separately for different tumor cohorts. One thousand

554 simulated datasets were analyzed with ActivePathways and those with at least one significantly

555 detected pathway counted towards false discovery rates. Additional simulations maintained the

556 positions of non-coding driver scores among gene scores and randomly reassigned protein-

557 coding driver scores, expectedly leading to a reduction in detected pathways as the input datasets

558 primarily included strong scores in protein-coding gene regions. Quantile-quantile analysis and

559 QQ-plots were used to compare p-value distributions of pathways discovered from true driver

560 scores, driver scores with shuffled driver scores, and driver scores shuffled entirely. To evaluate

561 robustness of ActivePathways, we randomly replaced a fraction of significant driver p-values in

562 input matrices (*P*<0.001) with insignificant p-values (*P*=1). We tested different fractions of missing

563 values (10%, 25%, 50%) across a thousand datasets of driver scores with randomly selected

564 missing data points and concluded that most cohorts included significantly enriched pathways

565 even with large fractions of missing data. To further evaluate robustness, we tested different

566 values of the Brown *P*-value threshold used to select the integrated gene list for pathway

567 enrichment analysis. The default parameter value ($P_{gene}$<0.1) was compared to alternative values

568 (0.001, 0.01, 0.05, 0.2). We concluded that ActivePathways found enriched pathways in most

569 tumor cohorts even at more stringent gene selection levels.

570 **Availability.** ActivePathways is freely available as an R package and source code on the GitHub

571 repository https://github.com/reimandlab/ActivePathways and the Comprehensive R Archive

572 Network (CRAN).

573

574 **Acknowledgements**

582

## 583  **References**

584  1    Creixell, P. *et al.* Pathway and network analysis of cancer genomes. *Nat Methods* **12**, 615-
585       621, doi:10.1038/nmeth.3440 (2015).

586  2    Reimand, J. *et al.* Pathway enrichment analysis of -omics data. *bioRxiv* **232835**,
587       doi:10.1101/232835 (2017).

588  3    Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nature*
589       *genetics* **45**, 1113-1120, doi:10.1038/ng.2764 (2013).

590  4    Hudson, T. J. *et al.* International network of cancer genome projects. *Nature* **464**, 993-
591       998, doi:10.1038/nature08987 (2010).

592  5    GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nature genetics* **45**,
593       580-585, doi:10.1038/ng.2653 (2013).

594  6    Subramanian, A. *et al.* Gene set enrichment analysis: a knowledge-based approach for
595       interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* **102**, 15545-
596       15550, doi:10.1073/pnas.0506580102 (2005).

597  7    Mi, H., Muruganujan, A. & Thomas, P. D. PANTHER in 2013: modeling the evolution of
598       gene function, and other gene attributes, in the context of phylogenetic trees. *Nucleic*
599       *Acids Res* **41**, D377-386, doi:10.1093/nar/gks1118 (2013).

600  8    Kaimal, V., Bardes, E. E., Tabar, S. C., Jegga, A. G. & Aronow, B. J. ToppCluster: a
601       multiple gene list feature analyzer for comparative enrichment clustering and network-
602       based dissection of biological systems. *Nucleic Acids Res* **38**, W96-102,
603       doi:10.1093/nar/gkq418 (2010).

604  9    Reimand, J. *et al.* g:Profiler-a web server for functional interpretation of gene lists (2016
605       update). *Nucleic Acids Res* **44**, W83-89, doi:10.1093/nar/gkw199 (2016).

606  10   Stein, L. D., Knoppers, B. M., Campbell, P., Getz, G. & Korbel, J. O. Data analysis: Create
607       a cloud commons. *Nature* **523**, 149-151, doi:10.1038/523149a (2015).

608  11   Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *bioRXiv* **162784** (2017).

609  12   Futreal, P. A. *et al.* A census of human cancer genes. *Nat Rev Cancer* **4**, 177-183,
610       doi:10.1038/nrc1299 (2004).

611  13   Huang, F. W. *et al.* Highly recurrent TERT promoter mutations in human melanoma.
612       *Science* **339**, 957-959, doi:10.1126/science.1229259 (2013).

613  14  Rheinbay, E. *et al.* Discovery and characterization of coding and non-coding driver
614      mutations in more than 2,500 whole cancer genomes. *BioRxiv* **237313** (2017).

615  15  Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene Ontology
616      Consortium. *Nature genetics* **25**, 25-29, doi:10.1038/75556 (2000).

617  16  Fabregat, A. *et al.* The Reactome Pathway Knowledgebase. *Nucleic Acids Res* **46**, D649-
618      D655, doi:10.1093/nar/gkx1132 (2018).

619  17  Brown, M. B. A Method for Combining Non-Independent, One-Sided Tests of Significance.
620      *Biometrics* **31**, 987-992 (1975).

621  18  Holm, S. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of*
622      *Statistics* **6 (2)**, 65–70 (1979).

623  19  Merico, D., Isserlin, R., Stueker, O., Emili, A. & Bader, G. D. Enrichment map: a network-
624      based method for gene-set enrichment visualization and interpretation. *PloS one* **5**,
625      e13984, doi:10.1371/journal.pone.0013984 (2010).

626  20  Reyna, M. A. *et al.* Pathway and network analysis of more than 2,500 whole cancer
627      genomes. *BioRxiv* **385294**, doi:https://doi.org/10.1101/385294 (2018).

628  21  Leiserson, M. D. *et al.* Pan-cancer network analysis identifies combinations of rare somatic
629      mutations across pathways and protein complexes. *Nature genetics* **47**, 106-114,
630      doi:10.1038/ng.3168 (2015).

631  22  Reyna, M. A., Leiserson, M. D. M. & Raphael, B. J. Identifying hierarchies of altered
632      subnetworks. *Bioinformatics* (2018).

633  23  Pulido-Tamayo, S., Weytjens, B., De Maeyer, D. & Marchal, K. SSA-ME Detection of
634      cancer driver genes using mutual exclusivity by small subnetwork analysis. *Sci Rep* **6**,
635      36257, doi:10.1038/srep36257 (2016).

636  24  Verbeke, L. P. *et al.* Pathway Relevance Ranking for Tumor Samples through Network-
637      Based Data Integration. *PLoS One* **10**, e0133503, doi:10.1371/journal.pone.0133503
638      (2015).

639  25  Sabarinathan, R., Mularoni, L., Deu-Pons, J., Gonzalez-Perez, A. & Lopez-Bigas, N.
640      Nucleotide excision repair is impaired by binding of transcription factors to DNA. *Nature*
641      **532**, 264-267, doi:10.1038/nature17661 (2016).

642  26  Curtis, C. *et al.* The genomic and transcriptomic architecture of 2,000 breast tumours
643      reveals novel subgroups. *Nature* **486**, 346-352, doi:10.1038/nature10983 (2012).

644  27  Quon, G. *et al.* Computational purification of individual tumor gene expression profiles
645      leads to significant improvements in prognostic prediction. *Genome Med* **5**, 29,
646      doi:10.1186/gm433 (2013).

647  28  Adams, S. *et al.* Prognostic value of tumor-infiltrating lymphocytes in triple-negative breast
648      cancers from two phase III randomized adjuvant breast cancer trials: ECOG 2197 and
649      ECOG 1199. *J Clin Oncol* **32**, 2959-2966, doi:10.1200/JCO.2013.55.0491 (2014).

650  29  Sabatier, R. *et al.* A gene expression signature identifies two prognostic subgroups of
651      basal breast cancer. *Breast Cancer Res Treat* **126**, 407-420, doi:10.1007/s10549-010-
652      0897-9 (2011).

653  30  Carpenter, R. L. & Lo, H. W. Regulation of Apoptosis by HER2 in Breast Cancer. *J*
654      *Carcinog Mutagen* **2013**, doi:10.4172/2157-2518.S7-003 (2013).

655  31   Wang, H. Y., Cheng, Z. & Malbon, C. C. Overexpression of mitogen-activated protein
656       kinase phosphatases MKP1, MKP2 in human breast cancer. *Cancer Lett* **191**, 229-237
657       (2003).

658  32   Denu, R. A. *et al.* Centrosome amplification induces high grade features and is prognostic
659       of worse outcomes in breast cancer. *BMC Cancer* **16**, 47, doi:10.1186/s12885-016-2083-
660       x (2016).

661  33   Belin, S. *et al.* Dysregulation of ribosome biogenesis and translational capacity is
662       associated with tumor progression of human breast cancer cells. *PLoS One* **4**, e7147,
663       doi:10.1371/journal.pone.0007147 (2009).

664  34   Guimaraes, J. C. & Zavolan, M. Patterns of ribosomal protein expression specify normal
665       and malignant human cells. *Genome Biol* **17**, 236, doi:10.1186/s13059-016-1104-z
666       (2016).

667  35   Varelas, X. The Hippo pathway effectors TAZ and YAP in development, homeostasis and
668       disease. *Development* **141**, 1614-1626, doi:10.1242/dev.102376 (2014).

669  36   Yu, F. X., Zhao, B. & Guan, K. L. Hippo Pathway in Organ Size Control, Tissue
670       Homeostasis, and Cancer. *Cell* **163**, 811-828, doi:10.1016/j.cell.2015.10.044 (2015).

671  37   Kolde, R., Laur, S., Adler, P. & Vilo, J. Robust rank aggregation for gene list integration
672       and meta-analysis. *Bioinformatics* **28**, 573-580, doi:10.1093/bioinformatics/btr709 (2012).

673  38   Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**,
674       415-421, doi:10.1038/nature12477 (2013).

675  39   Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-
676       associated genes. *Nature* **499**, 214-218, doi:10.1038/nature12213 (2013).

677  40   Wadi, L., Meyer, M., Weiser, J., Stein, L. D. & Reimand, J. Impact of outdated gene
678       annotations on pathway enrichment analysis. *Nat Methods* **13**, 705-706,
679       doi:10.1038/nmeth.3963 (2016).

680  41   Liberzon, A. *et al.* Molecular signatures database (MSigDB) 3.0. *Bioinformatics* **27**, 1739-
681       1740, doi:10.1093/bioinformatics/btr260 (2011).

682  42   Vaske, C. J. *et al.* Inference of patient-specific pathway activities from multi-dimensional
683       cancer genomics data using PARADIGM. *Bioinformatics* **26**, i237-245,
684       doi:10.1093/bioinformatics/btq182 (2010).

685  43   Warde-Farley, D. *et al.* The GeneMANIA prediction server: biological network integration
686       for gene prioritization and predicting gene function. *Nucleic Acids Res* **38**, W214-220,
687       doi:10.1093/nar/gkq537 (2010).

688  44   Parker, J. S. *et al.* Supervised risk predictor of breast cancer based on intrinsic subtypes.
689       *J Clin Oncol* **27**, 1160-1167, doi:10.1200/JCO.2008.18.1370 (2009).

690  45   Cline, M. S. *et al.* Integration of biological networks and gene expression data using
691       Cytoscape. *Nature protocols* **2**, 2366-2382, doi:10.1038/nprot.2007.324 (2007).

692